

Open Research Online

The Open University's repository of research publications and other research outputs

One document, many users: what happens when you re-purpose a document?

Conference or Workshop Item

How to cite:

King, David; Morse, David and Lyal, Chris (2013). One document, many users: what happens when you re-purpose a document? In: BioCuration 2013, 07-10 Apr 2013, Churchill College, Cambridge, UK.

For guidance on citations see [FAQs](#).

© 2013 ViBRANT

Version: Version of Record

Copyright and Moral Rights for the articles on this site are retained by the individual authors and/or other copyright owners. For more information on Open Research Online's data [policy](#) on reuse of materials please consult the policies page.

oro.open.ac.uk

One document, many users

1 Biologia Centrali-Americana

To assess the challenge around issues such as climate change and invasive species requires a baseline of historic data. We are fortunate in biodiversity that such data does exist in a rich body of literature. One such source of historic data is the Biologia Centrali-Americana (BCA), which documents the plant and animal life in Central America one hundred years ago, and which can be compared to contemporary species distributions. This valuable resource has recently been re-keyed and manually marked up by the INOTAXA project (<http://www.inotaxa.org/>) and is now being curated before wider release.

4 The taxonomist's view

```
<div type="taxon synonymy">
  <p elementid="BCA-aves-v3p1-2240">
    <hi rend="genus">
      <hi rend="italic">Vi reol ani us</hi>
    </hi>
    <hi rend="species">
      <hi rend="italic">mel itophrys</hi>
    </hi>,
    <bi bl rend="primary">
      <author>DuBus</author>,
      <title>Esq. Orn. </title>
```

The taxonomist needs to know the provenance of the taxon. Hence the mark-up is more than just the taxon name. In this example the taxon name is linked to the original describer of the taxon.

The taxonomist is also interested in typographical cues, such as the use of italic text.

Several information extraction tasks must be linked to provide a complete record.

7 Not marked up

A restricted range of entities is marked up. For example, Mexico and Guatemala are not recorded as countries.

8 What do computer scientists want?

Computer scientists prefer stand-off annotation so as to preserve the original text intact. This approach makes reuse of the text easier too.

The focus is on extracting chosen data. Frequently though this involves treating the text as independent tokens. Typographical cues are not considered, and collocation of terms is a specific task.

Note, computational linguists have a different view.

2 The re-purposing

Text mining has had some success in the recent, born-digital, bio-medical literature. Applying these approaches to the historic biodiversity literature is still in its infancy. One barrier is the lack of suitable corpora against which to develop and then test automated solutions. The ViBRANT project (<http://vbrant.eu/>) seeks to re-purpose the large volume of re-keyed data produced by INOTAXA to support the development of text mining solutions. However, this apparently straightforward task has thrown up many issues because biodiversity and computer scientists have different requirements of the mark up.

3 Additional challenges

This poster does not consider other challenges such as:

- Rekeyed data omits running headers, in the example below the re-keyed text omits VIREOLANIUS. 209
- OCR induced errors, in the example below the running header is identified as 'VIEEOLANIUS. 209', when it should read 'VIREOLANIUS. 209'; and the next line, which should read 'VIREOLANIUS.' is identified as 'VIKEOLANIUS.' Hence, we have two different incorrect recognitions.

5 Use or lose?

The genus name *Laniarius* is not marked up in the taxonomist's XML because it compares an African species to the Central American species being described. This work is concerned with documenting Central American species only.

For text mining purposes all taxonomic names are useful as training and testing data.

6 The computer scientist's view

T25 genus 1647 1658 Vi reol ani us
T26 speci fi cepi thet 1659 1670 mel itophrys

The computer scientist is concerned with one text mining operation at a time. The taxon name is not associated with an author, for example. These represent different name extraction challenges.

Contemporary computer science tools do not capture textual cues such as italics. Hence, potential semantic enhancements are more difficult to apply.

9 What do taxonomists want?

Taxonomists use inline XML to mark up their literature so that the enhancement and original literature are in the one document. There are three leading document level schemas:

- TaxonX, lightweight mark-up focused on taxon treatments (description of species).
- taXMLit, detailed mark-up focused on data curation, extraction and analysis.
- TaxPub, an extension of the National Library of Medicine DTD focused on layout and taxon names.

All three schemas have their advantages and shortcomings and can be used for different purposes.

See Penev L, Lyal CHC, Weitzman A, Morse D, King D, Sautter G, Georgiev T, Morris RA, Catapano T, Agosti D (2011) XML schemas and mark-up practices of taxonomic literature. In: Smith V, Penev L (Eds) e-Infrastructures for data publishing in biodiversity science. *ZooKeys* **150**: 89–116.

David King¹, David Morse¹, Chris Lyal²

1 Department of Computing, The Open University, UK.

2 The Natural History Museum, London, UK.

{David.King, David.Morse} @open.ac.uk
c.lyal@nhm.ac.uk



ViBRANT is funded by the European Union 7th Framework Programme within the Research Infrastructures group. Contract no. RI-261532. Period, Dec. 2010 to Nov. 2013. Coordinator: Dr Vince Smith. E-mail: enquiries@vbrant.eu

